

---

# **Machine Learning, Statistiques et Programmation**

*Version 0.1.407*

**Xavier Dupré**

oct. 27, 2018



---

## Table des matières

---

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Clustering</b>	<b>3</b>
2.1	k-means . . . . .	3
2.2	Mélange de lois normales . . . . .	18
2.3	Carte de Kohonen . . . . .	21
<b>3</b>	<b>Bases de Machine Learning</b>	<b>25</b>
3.1	Réseaux de neurones . . . . .	25
3.2	Classification à l'aide des plus proches voisins . . . . .	69
3.3	Liens entre factorisation de matrices, ACP, k-means . . . . .	77
3.4	Régression logistique, diagramme de Voronoï, k-Means . . . . .	82
<b>4</b>	<b>Natural Language Processing</b>	<b>111</b>
4.1	Complétion . . . . .	111
<b>5</b>	<b>Métriques</b>	<b>129</b>
5.1	Courbe ROC . . . . .	129
5.2	Confidence Interval and p-Value . . . . .	139
<b>6</b>	<b>Distances</b>	<b>157</b>
6.1	Distance d'édition . . . . .	157
<b>7</b>	<b>Graphes</b>	<b>163</b>
7.1	Distance between two graphs . . . . .	163
<b>8</b>	<b>Algorithmes</b>	<b>169</b>
8.1	Détection de segments . . . . .	169
<b>9</b>	<b>Pérégrinations d'un data scientist</b>	<b>179</b>
9.1	Répartir en base d'apprentissage et de test . . . . .	179
9.2	Corrélations non linéaires . . . . .	192
9.3	File d'attente, un petit exemple . . . . .	210
9.4	Optimisation avec données aléatoires . . . . .	216
9.5	Régression linéaire et résultats numériques . . . . .	221
9.6	Le gradient et le discret . . . . .	234
9.7	Régression quantile . . . . .	243

<b>10 API</b>	<b>251</b>
10.1 Machine Learning . . . . .	251
10.2 Traitement du langage naturel . . . . .	251
10.3 Source de données . . . . .	253
10.4 Graphes . . . . .	253
10.5 Image . . . . .	254
<b>11 Index</b>	<b>255</b>
11.1 Index . . . . .	255
<b>12 Galleries</b>	<b>299</b>
12.1 Le petit coin des data scientists . . . . .	299
12.2 Images . . . . .	304
12.3 Métriques . . . . .	304
12.4 Machine Learning . . . . .	311
12.5 NLP - Natural Language Processing . . . . .	328
<b>Bibliographie</b>	<b>357</b>

# CHAPITRE 1

---

## Introduction

---

J'ai commencé ce site un jour de nuit blanche, après avoir codé toute la journée et échangé avec mes collègues via des codes review. Les mails qui tombent sur mon portable, le énième commentaire sur la position des espaces dans mon code, une énorme angoisse m'étreignit alors. Les lignes de code défilaient comme les blocs de Tétris dans ma tête. Je me revois alors après une nuit blanche sur ce jeu, plus excité encore qu'après avoir du café à la fontaine. Je me suis dit qu'il fallait que j'arrête et que je remette un peu de sens dans tout ça.

Derrière le code, il y a les algorithmes et derrière encore des idées. Plein d'idées magnifiques, plein d'intuitions. Un code efficace n'a rien à voir avec sa lisibilité. C'est même souvent le contraire et il ne sera jamais lisible. Mais l'idée sous-jacente, elle, lorsqu'on l'a comprise, elle devient très claire. On n'a même plus besoin de l'écrire.

Il faudra probablement quelques années avant que ce site ne devienne conséquent, voire exhaustif. Il faut bien commencer quelque part. J'aurais pu écrire des pages wikipédia mais je n'aurais pas pu y mettre du code, des notebooks. J'ai commencé ce travail en latex lors de ma thèse, je me suis aperçu qu'il me restait quelques démonstrations à terminer.

Et en français... J'enseigne en français et puis l'offre est tellement riche en anglais.

### Installation

Les exemples de codes sont disponibles sous la forme d'un module python et des *notebooks* (page 299) accessibles sur le site.

```
pip install mlstatpy
```

Je n'ai pas la prétention d'être exhaustif. Vous la trouverez dans la documentation des bibliothèques de machine learning qui implémentent la plupart des algorithmes performants connus. Allez voir [scikit-learn](http://scikit-learn.org/stable/)<sup>1</sup> pour connaître ce qui marche. La documentation rend obsolète n'importe quel livre au bout de six mois.

---

1. <http://scikit-learn.org/stable/>



Ceci est un exemple, cliquez sur le lien de téléchargement pour obtenir le cours complet.

